JSC-07007

**NASA** NATIONAL AERONAUTICS AND SPACE ADMINISTRATION
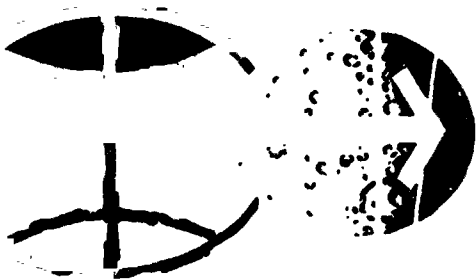
# JSC INTERNAL NOTE NO. 73-FM-53

## May 4, 1973

# THE APPLICABILITY AND EFFECTIVENESS

# OF CLUSTER ANALYSIS

## Mathematical Physics Branch

## MISSION PLANNING AND ANALYSIS DIVISION

## LYNDON B. JOHNSON SPACE CENTER
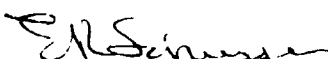### HOUSTON TEXAS

JSC-07907

JSC INTERNAL NOTE NO. 73-FM-53

# THE APPLICABILITY AND EFFECTIVENESS OF CLUSTER ANALYSIS

By D. S. Ingram, IBM, and
A. L. Actkinson, Mathematical Physics Branch

May 4, 1973

MISSION PLANNING AND ANALYSIS DIVISION

NATIONAL AERONAUTICS AND SPACE ADMINISTRATION

JOHNSON SPACE CENTER

HOUSTON, TEXAS

Approved: _____
Emil R. Schiesser, Acting Chief
Mathematical Physics Branch

Approved: _____
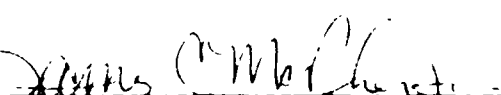John P. Mayer, Chief
Mission Planning and Analysis Division

CONTENTS

# FIGURES

THE APPLICABILITY AND EFFECTIVENESS OF CLUSTER ANALYSIS

By D. S. Ingram, IBM, and
A. L. Actkinson, Mathematical Physics Branch

## 1.0 SUMMARY

The objective of this internal note is to provide insight into the character-
istics which determine the performance of a clustering algorithm. It demonstrates
that, in order for the techniques which are examined to accurately cluster data,
two conditions must be simultaneously satisfied. The first condition is that the
data must have a particular structure, and the second is that the parameters chosen
for the clustering algorithm must be correct. By examining the structure of the
data from the Cl flight line, it is clear that there is no single set of parameters
that can be used to accurately cluster all the different crops. The effectiveness
of either a noniterative or iterative clustering algorithm to accurately cluster
data representative of the Cl flight line is questionable. This means that, in
order to use cluster analysis in its present form for applications like assisting
in the definition of field boundaries and evaluating the homogeneity of a field,
one must have extensive a priori knowledge. Modifications to existing techniques,
or entirely new techniques, are necessary for clustering to be a reliable tool for
representative data sets.

A modification to existing clustering methods is proposed. This involves
the use of goodness of fit tests to determine, in a quantitative manner, a measure
of the unimodality of a cluster. This also has applications to quantitatively
evaluating the homogeneity of test and training fields.

## 2.0 INTRODUCTION

Cluster analysis is a decision-making process in which similar measurements
are grouped together. The primary advantage of cluster analysis is that it is
not necessary to assume a statistical model for the data. Typical applications
which have been identified are evaluating field homogeneity, boundary definition,
selecting homogeneous data from nonhomogeneous data, and use as an unsupervised
classifier. An objective of this internal note is to determine the factors which
affect the ability of a clustering algorithm to perform these functions. These
factors are examined in view of the data analysis requirements associated with
processing multispectral scanner data for agricultural crops from the Cl flight
line.

For the clustering algorithms (ref. 1) which are examined to accurately cluster
data, two conditions must be simultaneously satisfied. First, the data must have a
particular structure, and, second, the correct parameters must be used in the clus-
tering algorithm. To demonstrate these conditions some experiments using two sets
of simulated data are described. The structure of the first set of data is such

2

that the clustering algorithm will accurately cluster the data. The statistics
of the second set of data are determined from the Cl flight line. These experi-
ments indicate that the results obtained by using existing cluster analysis tech-
niques to evaluate field homogeneity, boundary definition, selecting homogeneous
data from nonhomogeneous data, and as an unsupervised classifier are likely to have
little meaning unless one essentially knows the answer before the data are processed.
This is particularly significant because it means that it would be necessary to de-
termine the parameters to be used in the algorithm for each flight line and for
each different set of crops.

The key question which must be answered to make clustering a scientific tech-
nique rather than an art is whether a set of data is unimodal or multimodal. It
is proposed that two goodness of fit tests be investigated in order to quantify
the concepts of unimodality and homogeneity. This would be very valuable in de-
termining the appropriateness of the assumptions of the probability density func-
tion of the multivariate data. The assumption of the multivariate normal distri-
bution is used extensively in feature selection and pattern classification.


## 3.0 ANALYSIS


In this section clustering is initially described from an intuitive point of
view. The relationship between the form of the data and the result produced by a
clustering algorithm is investigated for some limiting cases. The results obtained
by processing data that corresponds to an agricultural image are presented for two
sets of simulated observations. The first set is chosen such that the algorithm
can produce accurate results. The statistics of the second set of data were de-
termined from the Cl flight line. Both a noniterative and an iterative algorithm
are used to process the data, which are representative of the Cl flight line.

As is demonstrated, the structure of the data corresponding to the agricultural
crops on the Cl flight line is such that no single set of parameters can be used to
accurately cluster the data. The cluster results is dependent on the parameters
used in the algorithm. Hence, any measure of field homogeneity is input-parameter-
dependent. The fundamental problem, then, is to determine whether or not a set of
data is unimodal.


### 3.1 Cluster Analysis

Cluster analysis is a decision-making process in which similar measurements
are grouped together. The performance of an algorithm to group data together de-
pends on the structure of the data. To illustrate this condition, consider two
sets of two dimensional data. In figure 1(a) the data are uniformly spaced in the
$x_1$, $x_2$ coordinates and in figure 1(b) the data are neatly grouped into three distinct
subsets.

In order to apply a cluster algorithm, a function must be used which determines
whether two observations are similar. For the sake of illustration let the similarity
function be a distance measure. If a measurement is within a specified radius of a
cluster mean, then that measurement is an element of the specified cluster. The
specified radius is a parameter of the algorithm. Consider the relationship be-
tween the number of clusters and the radius, R, for the data in figure 1(a). If
R is very small, then the number of clusters, N, equals the number of data points,
and if R is very large, then there is one cluster. As R changes from very small
to very large the structure of the N versus R cu⁻ would be similar to the graph

shown in figure 2(a) for the data of figure 1(a). The same procedure can be carried out for the data of figure 1(b). The limiting cases are the same; however, the number of clusters is constant over a range of R. Since the form of the data is known, it is clear that the correct number of clusters is three and that a value of R such that $R_1 < R < R_2$ is acceptable. This is equivalent to finding the set of parameters that produce the correct answer or to training the algorithm.

This example clearly shows that, in order for a clustering algorithm to effectively cluster data, two conditions must be simultaneously satisfied. The first condition is that the structure of the data must be such that the data can be clustered. If this condition is satisfied then it is necessary to choose the correct parameter in the algorithm. A value of R outside the range of $R_1 < R < R_2$ would not cluster the data correctly.

## 3.2 Simulated Data

The concept of clustering data from an image is further developed by considering two sets of simulated data. The clustering algorithms used include both a one-pass and an iterative technique. The spatial configuration of the classes in the image is similar to an agricultural scene. The first set of data is such that the clustering algorithm will effectively cluster the data. The second set of data is representative of the C1 flight line. Neither the noniterative nor the iterative algorithm is effective on the simulated C1 data.

3.2.1 Ideal case.- The simulated data generated for this case are described in reference 2. The noniterative clustering algorithm used is the CLUST1 option in ASTEP (ref. 1). Figure 3 illustrates the way the image is subdivided. For this example only two channels of data, 11 and 12 of reference 2, are processed. The mean ±1 standard deviation for each class are plotted in figure 4. Each element of the field is generated from a normal distribution with a mean of $\mu_i$ and standard deviation of $\sigma_i$ for $i = 1,5$ for each channel. The data are uncorrelated from channel to channel.

The data were clustered for several values of R with the condition that C = 2R. Although it is not obvious that the condition C = 2R will yield "best" results, this condition does appear to be a reasonable way of relating C and R. In each case the initial value for the maximum number of clusters was 20 and the initial values for the means of those clusters was 0. The results of the plot of the number of clusters versus R is shown in figure 5. The number of clusters is constant for values of $10 < R < 30$ and at each value of R the clusters are the same. The image map displayed by ASTEP is shown in figure 6 for R = 20. Each observation is classified correctly. This is exactly the result one would expect for the structure of the data in figure 4.

The question which one must ask is how to know that each of the five clusters is unimodal. For values of $32 < R < 50$ there are three clusters. The clusters are not the same three clusters for all values of R. However, for R = 32, 34, and 36 the three clusters are the same and one might suspect that there are three clusters instead of five.

4

3.2.2 <u>Simulated Cl data</u>.- The statistics of the crops along the Cl flight
line are listed in table 1. The mean ±1 standard deviation for each class in
channels 6, 10 and 12 is plotted in figures 7 and 8. The statistics listed in
table I were used to generate a data tape which represents the fields in figure 3.
The ellipses which are drawn in figures 7 and 8 have their principal axes parallel
to the measurement coordinates. This is not the case for Cl data, as the principal
axes of each of the ellipses would be inclined to the measurement coordinates. Ex-
amining figures 7 and 8 it is clear that the structure of the data is such it would
be difficult to find a set of parameters that could cluster all the crops. The
most obvious reason for this is the size of the standard deviation of Wheat2 as
compared to the difference between the means of the other crops, such as Alfalfa.

The noniterative algorithm described in reference 1 was used to process these
data with the same set of initial conditions described in section 3.2.1. The re-
sults of N versus R is shown in figure 9. There is clearly no well-defined
interval for which there are eight clusters. For this set of conditions the best
results appear to occur for R = 5 as shown in figure 10. It is possible to de-
cide what is best only because we know the answer.

Soybeans and Bare Soil are accurately classified. Corn1 and Oats are classi-
fied with fair accuracy. It is not possible to distinguish among Red Clover, Red
Clover2, Alfalfa, and Wheat 2. For this case the Red Clover field appears to be
nonhomogeneous while the Corn1 field appears to be homogeneous. As another illus-
tration of the concept of homogeneity consider figure 11. In this case R = 16
and the image divides nicely into two categories, B and C. The field labeled C
appears to be homogeneous and indeed is Wheat2. The field labeled B appears to
be homogeneous and indeed consists of seven different crops.

The same data were processed with ISODATA (ref. 4) using the same parameters
suggested in reference 5, namely DLMIN and STDMAX equal to 3.2 and 4.5, respectively.
For the best case (fig. 12), Red Clover and Alfalfa were indistinguishable and the
accuracy of the classification of Oats and Corn is fair. This case took 20 itera-
tions and used in NMIN value of 30. The term NMIN is the minimum number of points
allowed in a cluster. Changing the NMIN value to 15 resulted in Wheat, Red Clover,
Red Clover2, and Alfalfa being poorly classified while the accuracies of Corn and
Oat classification improved somewhat (fig. 13). The cases illustrate that the
choice of NMIN affects the accuracy of clustering. Changing NMIN may cause some ac-
curacies to improve while others deteriorate.

Chaining was applied in each of the above cases. The results without chaining
were much worse (figs. 14 and 15 for NMIN equal to 30 and 15, respectively).

Using a different channel set, channels 1, 6, 9, and 12, the ISODATA classifi-
cation maps were figures 16, 17, and 18, after 18, 19, and 20 iterations, respectively.
The value used for NMIN was 30, and no chaining was used. The 20 iterations case
(fig. 18) was less accurate than the corresponding three channel case given in
figure 14. Also, the results for iteration 18 were much better than those for itera-
tion 19.

These results demonstrate that the ISODATA classification is dependent on the
number of iterations, the number and choice of channels, and on the choice of NMIN.
No criteria currently exist for selecting these values without extensive a priori
knowledge. Even for the best choice, the accuracies were very poor for some crops.
The effectiveness of the iterative algorithm to cluster data representative of the
Cl flight line is questionable.

## 4.0 RECOMMENDATION

In order to determine whether a set of data consists of one or more clusters it is necessary to determine whether the data set is unimodal or multimodal. This can be determined in a quantitative manner by using goodness of fit tests. The two tests considered here are the classical chi-squared test and the Kolmogorov-Smirnov test. These could be applied to each potential cluster and the degree to which the cluster is unimodal could be determined.

A chi-squared random variable is defined as the sum of squares of independent standard normal variables (ref. 6). Let X be normally distributed with zero mean and unit variance; then the chi-squared random variable is

$$Q = X_1^2 + X_2^2 + \ldots + X_k^2 \tag{1}$$

where there are k independent values of X. The parameter k represents the number of degrees of freedom of the system. The chi-squared random variable can be related to the multivariate normal distribution by noticing that the quadratic form in the multivariate normal p.d.f. is a chi-squared random variable, that is,

$$p(x) = \frac{1}{(2\pi)^{n/2}\left|\sum\right|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T\sum{}^{-1}(x-\mu)\right\} \tag{2}$$

and

$$Q = (x-\mu)^T\sum{}^{-1}(x-\mu) \tag{3}$$

where x is the n × 1 random variable, μ is the mean vector and Σ is the n × n covariance matrix of the multivariate normal distribution. The chi-squared variable in equation (3) has n degrees of freedom.

The probability density function of a chi-squared random variable Q is

$$p(Q) = \frac{1}{2^{n/2}\Gamma(n/2)}Q^{(n-2)/2}e^{-Q/2} \quad Q > 0 \tag{4}$$

where n is the number of degrees of freedom and Γ is the gamma function. The cumulative distribution function (c.d.f.) of Q is

$$P(Q) = \int_0^Q p(Q)\,dQ \tag{5}$$

Equation (5) can be evaluated in closed form when n is even. The results for n = 2, 4, and 6 are

$$\underline{n = 2} \qquad P(Q) = 1 - e^{-Q/2} \qquad\qquad (6)$$

$$\underline{n = 4} \qquad P(Q) = 1 - e^{-Q/2}(1 + Q/2) \qquad\qquad (7)$$

$$\underline{n = 6} \qquad P(Q) = 1 - e^{-Q/2}\left[1 + \frac{Q}{2} + \frac{Q^2}{8}\right] \qquad\qquad (8)$$

The number of degrees of freedom is the same as the number of independent channels of a multispectral scanner.

If a data set has a multivariate normal distribution, then the numerically constructed p.d.f. and c.d.f. should match the functions generated by equation (4) and equation (5), respectively. The question of how well one function matches another introduces the concept of goodness of fit. Two goodness of fit techniques are developed, one related to the p.d.f. and the other related to the c.d.f. The advantages of the goodness of fit techniques is that it is possible to establish the percentile level of the fit.

In 1900 Pearson introduced the following measure (ref. 6), which is large when the differences $(f_{oi} - f_{ci})$ are large,

$$\chi^2 = \sum_{i=1}^{K} \frac{(f_{oi} - f_{ci})^2}{f_{ci}} \qquad\qquad (9)$$

where $f_{oi}$ is the ith observed frequency of occurrence, $f_{ci}$ is the ith expected or computed frequency, K is the number of measurements. It has been shown that $\chi^2$ is a chi-squared variable with k - 1 degrees of freedom. Hence, one could compute the frequency distribution of Q and evaluate $\chi^2$ for K intervals along the distribution and determine the percentile level from a table of percentiles of the chi-square distribution. In the case of an application to multispectral scanner data, the number of independent channels would determine the number of degrees of freedom to generate the p.d.f. given by equation (4), which is related to the computed frequency, $f_c$. Given the measurements the observed frequency could be constructed. Then K values along the $\chi^2$ axis could be chosen and equation (9) evaluated. The use of a table of percentiles of the chi-square distribution would determine the accuracy to which the data base follows the chi-squared assumption.

A method of determining the goodness of fit based on the distribution function uses the Kolmogorov-Smirnov statistic. If $P(x)$ is the theoretically constructed cumulative distribution function and $P_0(x)$ is the numerically constructed cumulative distribution function then the Kolmogorov-Smirnov statistic is

$$D = \max_{\text{all } x} |P(x) - P_0(x)| \qquad (10)$$

This situation is illustrated in figure 19. The value of D and the number of samples determined the accuracy to which $P_c(x)$ approximates $P(x)$. In the case of multispectral scanner data the number of independent channels determines the value of n to be used in equation (5). $P_c(x)$ would be generated from the experimental data. The value of equation (10) and the number of samples would be used as inputs to a table of acceptance limits for the Kolmogorov-Smirnov test of goodness of fit.

The effectiveness of the chi-squared statistic and the goodness of fit tests should be evaluated by using synthetic data. This is an effective procedure for checking out the implementation and pc·er of the algorithms. Synthetic data which are representative of aircraft and spacecraft data should be generated and analyzed. This will provide insight into the applicability of the statistical tests for different data bases.

Actual remotely sensed data from aircraft and spacecraft should be processed to obtain better insight into the characteristics of the data base. The topics of field homogeneity, feature selection, pattern classification, and error analysis should be investigated in terms of the characteristics of the data base.

## 5.0 CONCLUSION

This internal note has demonstrated that, in order for the techniques which are examined to accurately cluster data, two conditions must be s·multaneously satisfied. The first condition is that the data must have a particular structure, and the second is that the parameters chosen for the clustering algorithm. must be correct. By examining the structure of the data from the C1 flight line, it is clear that there is no single set of parameters that can be used to accurately cluster all the different crops. The effectiveness of either a one-pass c iterative clustering algorithm to accurately cluster data representative of the C1 flight line is questionable. This means that, in order to use cluster analysis in its present form for applications like assisting in the definition of field boundaries and evaluating the homogeneity of a field, one must have extensive a priori knowledge.

Modifications to existing techniques, or entirely new techniques, are necessary for clustering to be a reliable tool for representative data sets. This involves the use of goodness of fit tests to determine, in a quantitative manner, a measure of the unimodality of a cluster. This also has applications to quantitatively evaluating the homogeneity of test and training fields.

TABLE I.— MEAN/STANDARD DEVIATION OF C1 AGRICULTURAL DATA

Channel Number

| Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SOY | 84.46 / 2.40 | 79.76 / 2.58 | 61.06 / 1.82 | 62.32 / 1.81 | 85.78 / 3.55 | 87.92 / 3.37 | 63.90 / 2.17 | 84.19 / 3.78 | 70.55 / 3.20 | 81.87 / 3.47 | 91.57 / 5.83 | 73.28 / 3.60 |
| CORN 1 | 83.89 / 3.02 | 77.89 / 3.00 | 59.34 / 1.78 | 59.95 / 1.83 | 82.33 / 3.10 | 86.75 / 2.49 | 61.93 / 1.82 | 76.82 / 3.18 | 61.93 / 2.97 | 72.61 / 3.06 | 105.94 / 7.37 | 81.56 / 4.58 |
| CORN 2 | 75.90 / 1.59 | 71.53 / 1.77 | 56.07 / 1.30 | 56.48 / 1.22 | 75.55 / 1.96 | 81.21 / 1.72 | 59.51 / 1.46 | 72.36 / 2.25 | 59.10 / 2.12 | 71.66 / 2.30 | 114.78 / 6.31 | 88.73 / 3.80 |
| OATS | 75.78 / 3.00 | 73.25 / 2.88 | 57.15 / 1.93 | 58.29 / 1.94 | 79.60 / 4.01 | 84.63 / 3.85 | 63.07 / 2.39 | 85.15 / 4.35 | 71.33 / 4.34 | 87.73 / 4.21 | 106.97 / 7.00 | 85.69 / 4.97 |
| WHEAT 1 | 72.16 / 2.82 | 71.40 / 3.00 | 57.47 / 1.98 | 59.64 / 2.03 | 80.47 / 3.91 | 81.15 / 3.85 | 62.82 / 2.68 | 92.74 / 5.23 | 83.75 / 5.26 | 95.68 / 7.33 | 81.18 / 7.27 | 69.35 / 5.55 |
| WHEAT 2 | 80.42 / 2.69 | 80.86 / 4.50 | 64.42 / 3.27 | 67.35 / 3.85 | 99.32 / 8.51 | 102.06 / 9.07 | 76.86 / 6.46 | 116.76 / 11.79 | 100.79 / 9.64 | 116.27 / 12.76 | 98.31 / 10.97 | 77.83 / 7.36 |
| RED CLOV | 72.34 / 2.02 | 69.38 / 2.06 | 54.00 / 1.31 | 54.18 / 1.44 | 72.19 / 2.77 | 79.96 / 3.18 | 57.60 / 1.90 | 68.92 / 2.79 | 55.35 / 2.40 | 76.58 / 3.37 | 140.02 / 14.54 | 107.96 / 8.87 |
| RED 2 | 71.53 / 1.57 | 68.36 / 1.66 | 53.24 / 1.31 | 53.71 / 1.21 | 71.28 / 1.80 | 78.54 / 1.99 | 56.92 / 1.22 | 67.68 / 1.88 | 54.17 / 1.86 | 71.15 / 2.77 | 121.77 / 14.29 | 94.59 / 8.79 |
| ALFALFA | 76.84 / 1.88 | 71.68 / 2.37 | 55.27 / 1.71 | 55.47 / 1.87 | 75.13 / 3.16 | 84.59 / 2.17 | 60.08 / 1.81 | 69.64 / 4.14 | 54.42 / 3.80 | 77.97 / 2.41 | 154.21 / 14.25 | 114.00 / 8.28 |
| RYE | 80.12 / 2.19 | 79.79 / 2.56 | 63.68 / 1.83 | 65.36 / 1.78 | 93.92 / 2.95 | 96.71 / 2.37 | 72.57 / 1.93 | 102.91 / 3.80 | 86.13 / 3.90 | 99.84 / 3.86 | 96.05 / 4.18 | 75.86 / 2.48 |
| SOIL | 90.05 / 1.80 | 85.31 / 1.82 | 65.42 / 1.35 | 66.84 / 1.17 | 92.52 / 1.78 | 89.47 / 1.50 | 67.16 / 1.24 | 95.12 / 1.82 | 83.46 / 1.47 | 90.24 / 1.90 | 71.81 / 2.79 | 60.38 / 1.76 |

Figure 1.- Uniform and grouped data.

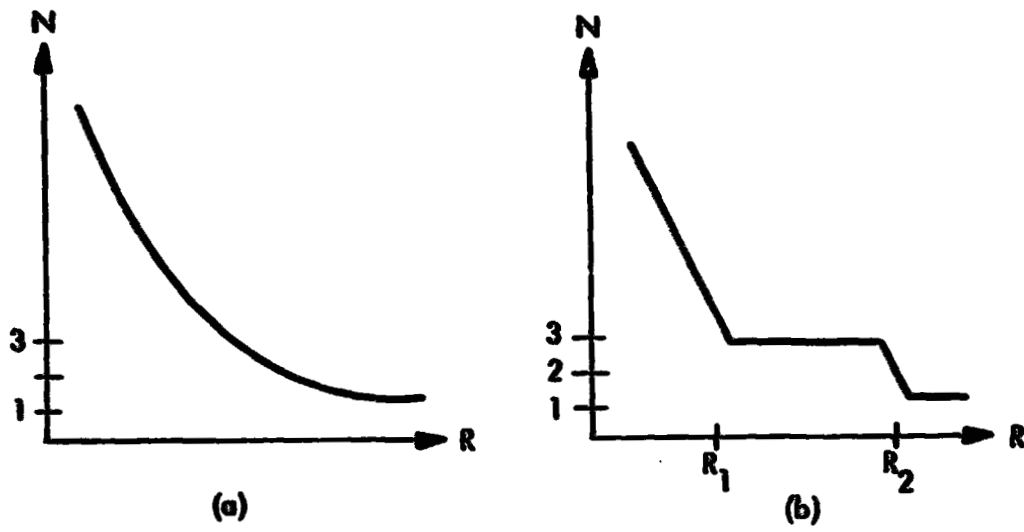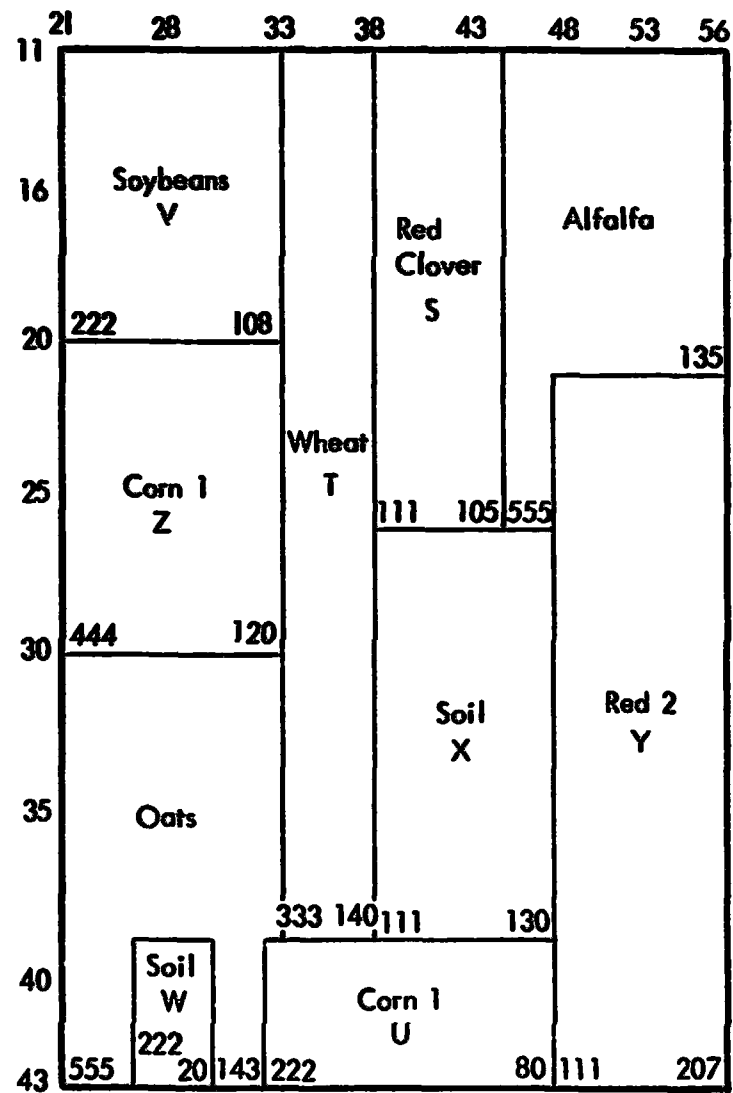Figure 2.- N versus R for the data of figure 1.

For each field:

Letter in middle corresponds to FIELD ID
Number in lower left corresponds to Class ID
Number in lower right corresponds to number of pixels
The agricultural crop is used in simulating C 1 data.

Figure 3.- Data image (fig. 1 of ref. 2).

12

Figure 4.- Simulated data from the SERID program.

Figure 5.- number of clusters versus R for the ideal case.

COMPLETE IMAGE

```
0000000001111111111222222222233333333
123456 8901234567890123456789 0123456
```

```
 1  8883988888983UUDUULEEELEELCCCCCCCLCCCC
 2  888388888833 UUDJULEEEEEELCLCCCLCC CCCC
 3  8898398838383UUDJULEEEEEELCCCCCCCCCCCCC
 4  8988388888333UUDJULEEELEELCCCCCCCCCCCCC
 5  8883888838333UUDJULEEEEEELCCCCCCCCCCCCC
 6  8888388888388UUDJULEEELEELCCCCCCCCCCCCC
 7  3383888888338UUDJULEEELEELCCCCCCCCCCCCC
 8  338338893888UUDJULEEELEELCCCCCCCCCCCCC
 9  88333888888 UUDJULEEEEEELCCCCCCCCCCCCC
10  AAAAAAAAAAAA UUDJULEEELEFLCCCCCCCCCCCCC
11  AAAAAAAAAAA UUDJULEEEELELLCCCCCEEEEEEEE
12  AAAAAAAAAAA UUDJULEEELELLLCCEEEEEEEEE
13  AAAAAAAAAAA UUDJULEEEELFLLCCEEEEEEEE
14  AAAAAAAAAAA UUDJULEEEEELLCCCLELEEEEE
15  AAAAAAAAAAA UUDJULEEELELLCCCLELEEEEE
16  AAAAAAAAAAA UUDJULEEELELEELEEEEEEE
17  AAAAAAAAAAA UUDJULEEEEELLEELEEEEEEE
18  AAAAAAAAAAA UUDJULEEELEELEELEEEEEEE
19  AAAAAAAAAAA UUDJULEEELEELEELEEEEEEE
20  CCCCCCCCCCC UUDJULEEELELLEELEEEEEE
21  CCCCCCCCCCC UUDJULEEELELEELEEEEEEE
22  CCCCCCCCCCC UUDJULEEELELEELEEEEEEE
23  CCCCCCCCCCC UUDJULEEELELEELEEEEEEE
24  CCCCCCCCCCC UUDJULEEELELEELEEEEEEE
25  CCCCCCCCCCC UUDJULEEEELELEELEEEEEE
26  CCCCCCCCCCC JJJULEEELELEELEEEEEE
27  CCCCCCCCCCC JJDJULEELEELEELEEEEEEE
28  CCCC CCCCCC UUDJULEEELELEELEEEEEE
29  CCCC8888 CCC 88888888888 EELEEEEE
30  CCCC8888 CCC 888888888888 EELEEEE
31  CCCC8888 CCC 888888888888 EELEEEE
32  CCCC8888 CCC 888888888888 EEEEEEE
33  CCCC8838 CCC 888888888888 EELEEEE
```

Figure 6.- Image map from the noniterative clustering technique
for the ideal case, R = 20.

Figure 7.- Intensity of channels 6 and 10 for Cl data.

16



Figure 8.- Intensity - channels 6 and 12 for C1 data.

Figure 9.- N versus R for simulated Cl data.

COMPLETE IMAGE

```
                    00060000011111111112222222222233333333
                    123456789C123456789D1234567690123456

 1    BBBBUBBBCBBBPGAAAAIAIIANNNANINANNNNN
 2    BCBBBBBBBCBBAMAAMAIAIIINIIINAINNIANA
 3    CBBBBBBBBBBBPLAAJAIAIIIAANIMAININNNN
 4    BBBCBBBBBBBBKAAAAAIAAIAANNMMANAIANNI
 5    BEBBBCBBBBBBLAMAAAIAAAAAIINAAIIINANN
 6    BBBBBBBBBBBBMAQAMAIIAIAANNAIAAINNAAI
 7    BBBEBBBBBBBBMADLAAIAAIAIMNANAINAIINN
 8    CBBEBBBBBBBBUAMAAIIAAIAIIINAAINNNNNN
 9    BBBBBBBBBBBBMAAAAIAAAAANIIIAANNIMNAA
10    CCCCCCCCCLCCAWAAAAIAAAAIMAAAAAIIANI
11    CCCCCCCCCLCCPMAAAAIAIAAMAISASSSWASI
12    CCCCLCCCCCCCAAAKUAAIIIAAAIMSASSAAIAS
13    CCCCCCCCCCCCPAAALAIIAAAIMIAAACASISAS
14    CCCCCCCCCLCCPAAAJIAAAAAIIAIAASASSIAS
15    CCCCCCCCCCBCMAAAAIIIIAIMIAISAIISASSA
16    CCCCCCCCCCUCCPAAJAPLLEELEEPSASISIASS
17    CCCCCCBBCLCCAAJMALEEEELEEEEAAFCAIAIS
18    CCCBCCCCCCCCAAAAMEEEELEEEELAMASSSASS
19    CCCCCCBCCCCCAAAAMLLEELLEEEEAASSSSSAA
20    UUULBUUDDUUUDAAAPAELELLELELLAASASAASA
21    ULULDUUDDUUUFAAAALLLEELELEESAAASASSA
22    UUUUDDDDCFUUFAAAAFELEELELEESSASAACSI
23    UUULUBDDDUGDAAAAALEEEELELELASASSSAAS
24    UUUUDDDDUUUUMAPAALLLEELELELAAAASIAIC
25    CUUUUDDODLUCAMJMALELELELEEASISSSCS5
26    UUUUBBDDUUUDAAALAELLELLLLLLCSSSSSSAC
27    UUUUDODDUUBUPPRARELEEEELELCAASSSSSC
28    UUBUBUDDDUUULJARAALLLEELELLLAASSSSSAS
29    DUBBEEEEUUUCCCCACCCCCACCLCCAASAAASAS
30    UUUUEEEEUUUCLCAACCCACCCACCLAIIAAASAA
31    UUULEEEEUUUCCHMCCCCCCCCCLACHISSSAASS
32    UUULEEEEUUUCCCACAACLCCCCCAAASACSSSSS
33    UUUUEEEEUUUCCACCACAACCCCSACSSASSISST
```
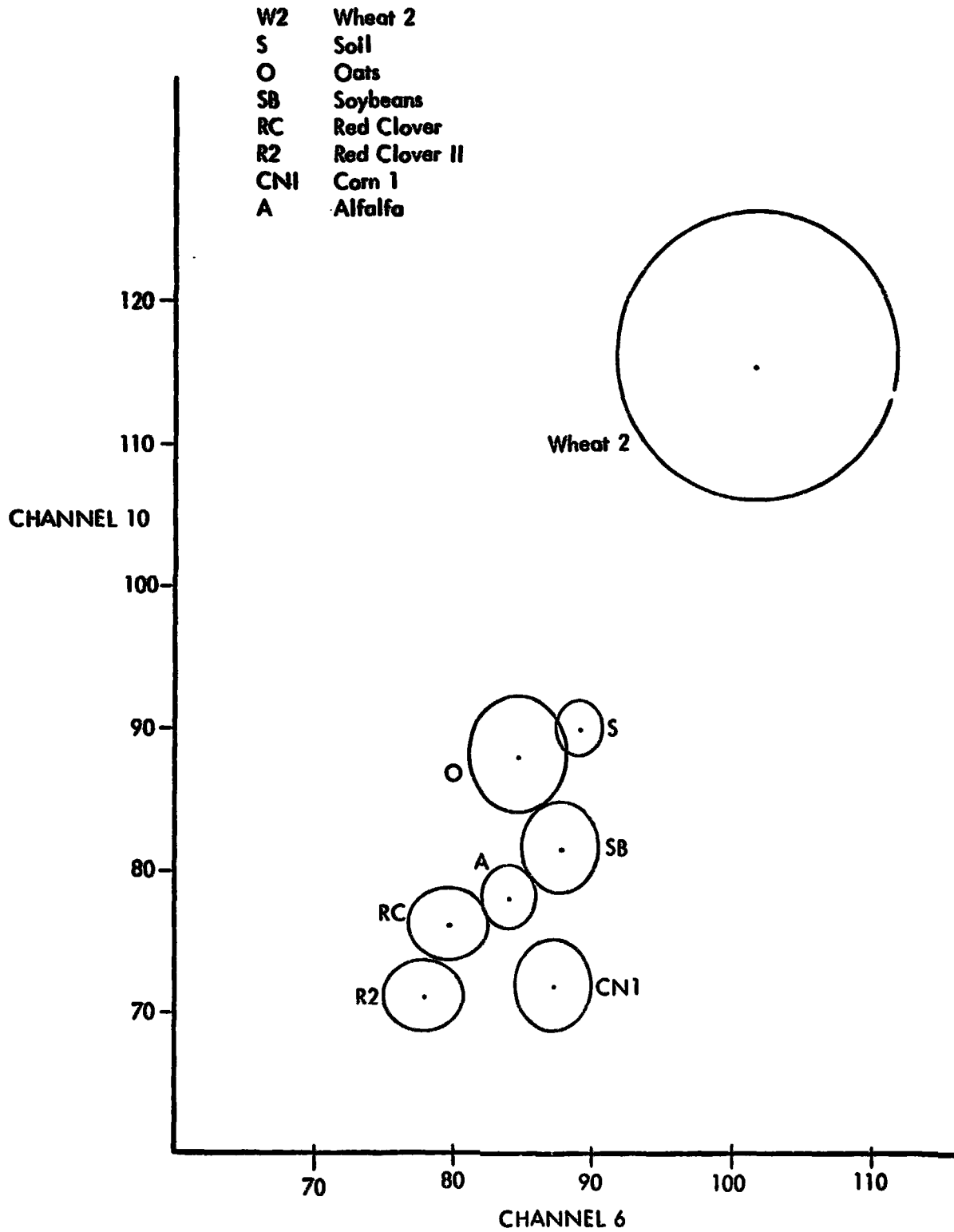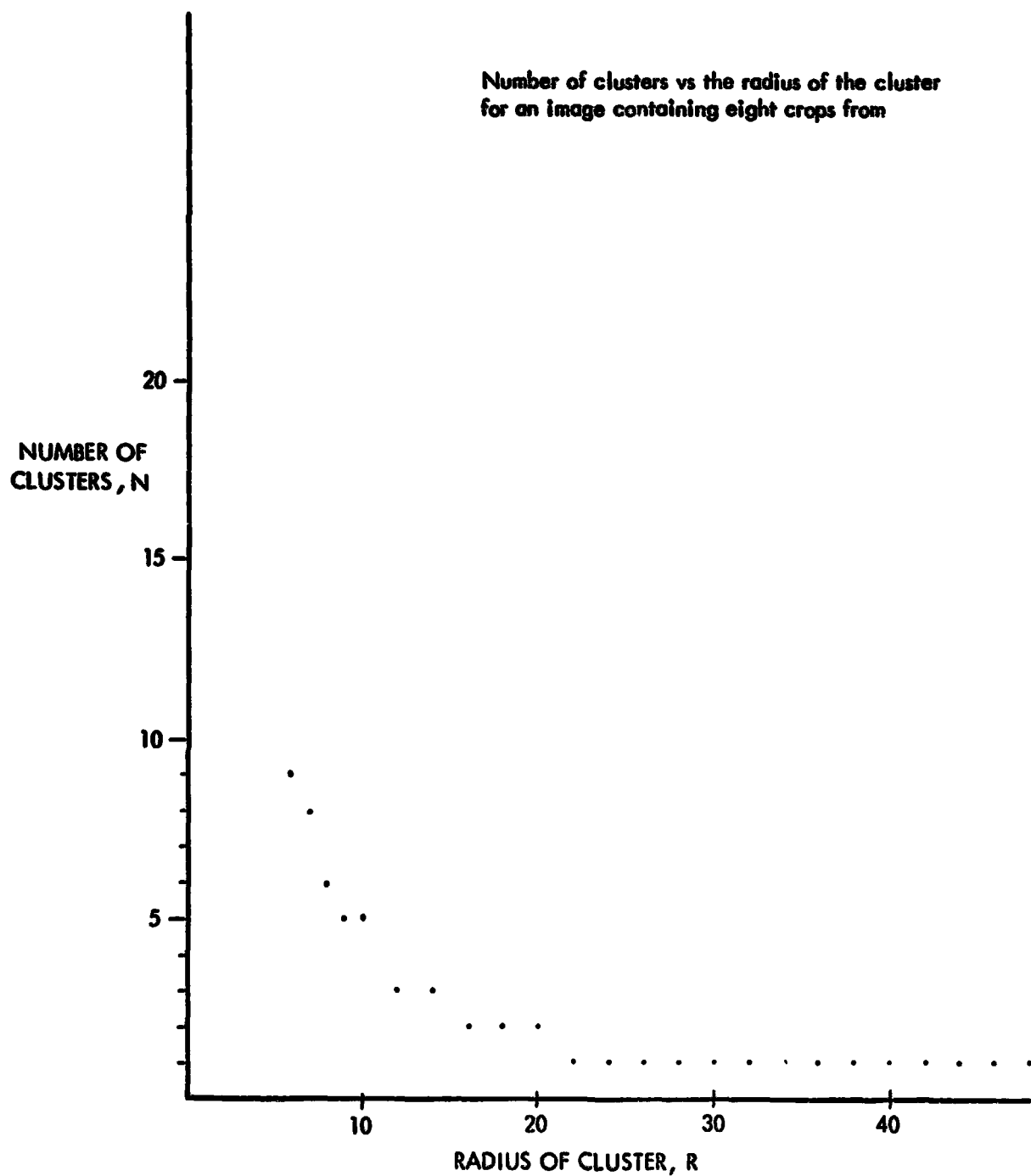
Figure 10.- Image map from the noniterative clustering technique for R = 5.

Figure 11.- Image map from the noniterative clustering technique for R = 18.

Figure 12.- ISODATA results (channels 6, 10, 12;
NMIN = 30; no. of iterations = 20; with chaining.

Figure 13.- ISODATA results (channels 6, 10, 12;
NMIN = 15; no. of iterations = 20; with chaining

Figure 14.- ISODATA results (channels 6, 10, 12; NMIN = 30; no. of iterations = 20; no chaining

```
00000000000000000000000O0000O000000000000
0000000000111111111111222222222223333333
1234567890123456789012345678901234456
```

```
 1   ./.,//,/.///II+IEJW-WWW--=W-W-===-==
 2   /A././///.//+KFIKW==WWJ-W=W-JW--W==-
 3   ./.,/...///...I+@++SWWW=J-=-W--W-=---=
 4   ////....///./GIKIGWW=S=J=----==-==-=W
 5   /./.,/A.//./+.+EG+=-W--=WW-JWWWW=W-=
 6   //..,//./////.K+.BIW=W=W==-===W-====
 7   /./////////,E+KE=-J=JWWW=-=-WW-W=W==
 8   ./.,//,/.,,///+BEK+WWWSW=WWW===W=-==-=
 9   //.,/.///.///E.K+EWJJJJW--WWJ=---W-==J
10   AACAAAAAAAAAK=EGE-WWJJ-=-WW==-W==W==
11   AAAAAAJDAAAAIK+G==W=J=--=WSSSCC=JSW
12   AAAAAAAAACAAAGIGIK===WW===W-SSS$AAWWC
13   AAAAAAAAAAAAAGEGE=SWW=WJW=W=WJCCSSCCS
14   ACAAACAAAAAAGIKEE=CW=W=WWW=CCS$SJSSC
15   AAAAAAAAAA/CEI+GIWWWWW==I=-SCWSCJSSC
16   AAAAC.AAA/AAIIGE+BBBBBBBBBBSSS$SSACC
17   AAAAAA//AAAAII+EEBBBBBBBBBBCWCASSSS
18   AAA/AAAAA/CAKG+FIBBBBBBBBBBCSWSSSC$J
19   AAAAAA/AAAAA+GEEBBBBBBBBBBCCJSSCSCC
20   .DFF.DHDFFD+GGG+BBBBBBBBBKSJCSCCCC
21   @HFFDDHFFHD.GI+GGBBBBBBBBBCCSCSCCJC
22   FDDFFH@F.FDDK++E+BBBBBBBBBSSCCCSAJS
23   DDDHD.HF@HHD@KK=+BBBBBBBBBSCS$SCCCC
24   HFFF@FFFDH@DEGIIGBBBBBBBBBCWCCSWCWC
25   ADD@F@H@FHD.E++K+BBBBBBBBBWS=JSSAS$
26   .HDD,.HFDDFDEKI+FBBBBBBBBBCSSSSSSAC
27   DDHHD@@F.J.DGEIE+BBBBBBBBBCSSCJCJSC
28   DH.H.@FFDFDD+IIGKBBBBBBBBBSCSCCC=C
29   DD..BABBB.I-AAAAAAAAAAAAAAAJSSWCSCJ
30   .FFFBBBB@@DAAAAAAAAAAAAAAAASSASSSCS
31   HDF.BBBBD@FCA/AAAAAAAAA/AJACSSCCCSJC
32   D.H.BBBBFFDAAAAAAAAAAAAAAA/BSJCCSCSC
33   D.H@BBBB@FFAAAAAACAAAAAASAABSA$SWJJC
```

Figure 15.- ISODATA results (channels 6, 10, 12;
NMIN = 15; no. of iterations = 20; no chaining

24

Figure 16.- ISODATA results chart NMIN = 30; no. of iterations =

Figure 17.- ISODATA results (channels 1, 6, 9, 12; NMIN = 30; no. of iterations = 19; no chaining

```
000000000000000000000000000000000000000000
000000000011111111112222222222233333333
1234567890123456789012345678901234566
```

```
 1  CCCCCCCC/CAC+BBB+B=-S==+=--=-+++=-++
 2  C/CCCCCCC/#C+BB/+-=====S=+==--S=--=++=
 3  /CCCC/CCCCCC++B+-B=-S+S==--=--=--=---+
 4  CCCCCCCCCACC++B+-B==-S+S+----==+++=-==
 5  CACCCACCCCCCB++++B+---=--+==-B==-===--+
 6  CC/CCCCCCCCC+BB+B====-=-+--+++++=--+++=
 7  CCCCCCCCCCCCC+BB+B+S+S==++++--==--=++++
 8  ACCCCC/C/CCCB++++BS=SS=+S==++++=++---+-
 9  CCA#CCCCCCCCBB#++SSSBS==+==S+--=--++,
10  777C77C7CC70+B+B+-=SSS-+B==+===--==+++
11  /////C///CC//B++++=-B+--+++SSSBB+SS=
12  //AC//C/////CB++B+=++++=====-+SSSSBB==B
13  //C//C////////+++BBCSS===S+=+==SBBSSBBS
14  /////////C///+++++B-B=-S++==+BBSSSSSSB
15  ////C/C/CACB+BB+SS===-=+=BB==BBSSB
16  /C////////CC/BBBB+#######################SS=SSBBB
17  //C/C/CC/C//B++BB#######################=BBSSSSS
18  C//CA//C/C//B+++B#######################S=SSSBSS
19  C/////C/C/C//C++B+B#######################BBSSSBSBB
20  AAAAA++++AAAA+++BB#######################SSBSBBBB
21  ++AAAA+A+++A+B+++##########################BSBSBBSB
22  A+AAA+++/AB+BB+++##########################SBBBSBS=
23  +A++AA+A+++A+B+BBC#######################SBSSSBBBB
24  +A+/+++/A+/ABBB++#######################SBBS=B=B
25  AA//A+++A+AABB+BB#######################S+SSSBSS
26  C+++AA+A+//+BB+BB#######################SSSSSSBB
27  +A++/++/A+A++/+B+#######################SSBBBSSB
28  ++A+A/AAAAAAB+BB+#######################SBSBBB+B
29  A+AA####+A+///CB////////////C//BSSS BSBS
30  A++A####+/A///////CC//C//////BSSBSSSBS
31  +AAA#####+A/CC/C///////Ca/CBSSBBBSSB
32  +C+A####+AA+//////////A/A////CBSSBBSBSB
33  +A++####+AA+C//C//CC////C/BSBSS=SSB
```

Figure 18.- ISODATA results (channels 1, 6, 9, 12;
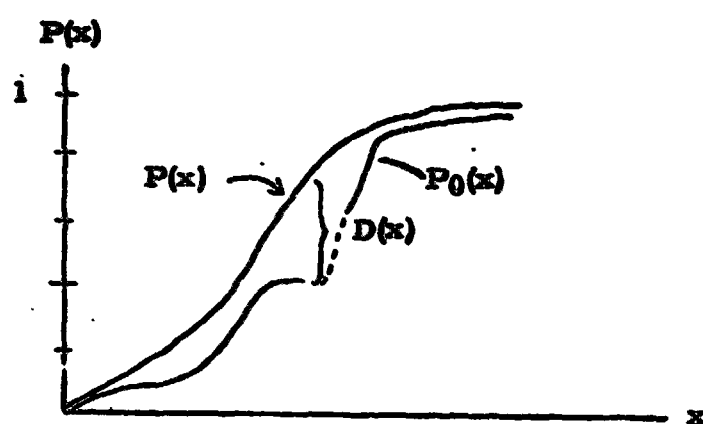NMIN = 30; no. of iterations = 20; no chaining

Figure 19.- Kolmogorov-Smirnov statistic.

## REFERENCES

1. Detchmendy, D.; and Wylie, A.: RTCC Requirements for Implementation of a
   Clustering Program in the ERIPS. MSC IN 72-FM-268, Jan. 18, 1973.

2. Wylie, A. D.: Simulated Test Case for Verifying the SERID Program.
   MSC Memo FM8(72-154), June 21, 1972.

3. Detchmendy, D., et al: User's Guide and Software Documentation for the
   Algorithm Simulation Test and Evaluation Program (ASTEP). MSC IN 72-FM-19C,
   Aug. 4, 1972

4. Kan, E. P. F.; and Holley, W. A.: More on Clustering Techniques with Final
   Recommendations on ISODATA. Lockheed Electronics Co., HASD, Houston, Texas,
   Technical Report 640-TR-112, May 1972.

5. Kan, E. P. F.: The Latest Version of ISODAT(A)/ISOC. Lockheed Electronics Co.,
   HASD, Houston, Texas, Technical Memo 63-0257-2115, Sept. 1972.

6. Lindgreen, B.: Statistical Theory. The MacMillan Co., New York, 1962.